

# Notes on Matrix Derivative in Neural Network

Wu Xiaojian

August 3rd, 2019

## 1 Notation

Let lower case letters  $a, b, c, \dots$  be scalar values, bold lower case letters  $\mathbf{v}, \mathbf{x}, \mathbf{y}, \mathbf{z}, \dots$  be vectors and bold upper case letters  $\mathbf{A}, \mathbf{B}, \mathbf{C}, \dots$  be matrices. Detailedly, we have

$$\begin{aligned} a, b, c, \dots &\in \mathbb{R} \\ \mathbf{v} &= [v_1, v_2, \dots, v_n]^T \quad \text{for } \mathbf{v} \in \mathbb{R}^n \\ \mathbf{A} &= \begin{bmatrix} a_{11} & \cdots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{m1} & \cdots & a_{mn} \end{bmatrix} \quad \text{for } \mathbf{A} \in \mathbb{R}^{m \times n} \end{aligned}$$

We introduce a notation  $\delta_{ij,pq}$ , which will be used in the formulas later, as below:

$$\delta_{ij,pq} = \begin{cases} 1 & \text{if } i = p \text{ and } j = q \\ 0 & \text{otherwise} \end{cases}$$

in some cases, we are to leave out one term and write  $\delta_{i,p}$ , meaning 1 if  $i = p$ , and 0 otherwise.

## 2 Derivative of Scalar, Vector and Matrix

Let  $L \in \mathbb{R}$ ,  $\mathbf{x} \in \mathbb{R}^n$ ,  $\mathbf{y} \in \mathbb{R}^m$ ,  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , we define derivative operations among them:

- scalar value and vector

$$\frac{\partial L}{\partial \mathbf{x}} = \left[ \frac{\partial L}{\partial x_1}, \frac{\partial L}{\partial x_2}, \dots, \frac{\partial L}{\partial x_n} \right] \quad (1a)$$

$$\frac{\partial \mathbf{x}}{\partial L} = \left[ \frac{\partial x_1}{\partial L}, \frac{\partial x_2}{\partial L}, \dots, \frac{\partial x_n}{\partial L} \right]^T \quad (1b)$$

**NOTICE:** The Derivative of  $L$  to  $\mathbf{x}$  is a **row vector** of  $1 \times n$  but not a column vector of the same shape of  $\mathbf{x}$ , while the derivative of  $\mathbf{x}$  to  $L$  is a column vector with the shape consistent to  $\mathbf{x}$ .

- vector to vector

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y_1}{\partial x} \\ \frac{\partial y_2}{\partial x} \\ \vdots \\ \frac{\partial y_m}{\partial x} \end{bmatrix} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \cdots & \frac{\partial y_1}{\partial x_n} \\ \frac{\partial y_2}{\partial x_1} & \frac{\partial y_2}{\partial x_2} & \cdots & \frac{\partial y_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \frac{\partial y_m}{\partial x_2} & \cdots & \frac{\partial y_m}{\partial x_n} \end{bmatrix} \quad (2)$$

we give a name to such a matrix as **Jacobian Matrix**.

- scalar value to matrix

$$\frac{\partial L}{\partial \mathbf{A}} = \begin{bmatrix} \frac{\partial L}{\partial a_{11}} & \frac{\partial L}{\partial a_{21}} & \cdots & \frac{\partial L}{\partial a_{n1}} \\ \frac{\partial L}{\partial a_{12}} & \frac{\partial L}{\partial a_{22}} & \cdots & \frac{\partial L}{\partial a_{n2}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial L}{\partial a_{1n}} & \frac{\partial L}{\partial a_{2n}} & \cdots & \frac{\partial L}{\partial a_{nn}} \end{bmatrix} \quad (3)$$

Like EQUATION(1a), the result matrix is obtained by applying derivative of  $L$  to each entry of  $\mathbf{A}$  and taking **transpose operation** on it, so that the shape is  $n \times m$  but not  $m \times n$ .

In multi-variable calculus, we know that if  $L = L(y_1, y_2, \dots, y_n)$  and  $y_i = y_i(x_1, x_2, \dots, x_m)$  ( $i = 1, 2, \dots, n$ ), then

$$\frac{\partial L}{\partial x_k} = \sum_{i=1}^n \frac{\partial L}{\partial y_i} \frac{\partial y_i}{\partial x_k}$$

(**NOTICE:** In single-variable derivative, say,  $z = f(y)$ ,  $y = g(x)$ , by **chain rule** we have  $\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$ . But in partial derivative, we cannot conclude that

$$\frac{\partial L}{\partial y_i} \frac{\partial y_i}{\partial x_k} = \frac{\partial L}{\partial x_k})$$

Let's regard these functions in another way. Since  $L$  takes  $n$  arguments as input, it's natural to rewrite  $L$  as  $L = L(\mathbf{y}) = L([y_1, y_2, \dots, y_n])$ , where  $y_i = y_i(\mathbf{x}) = y_i([x_1, x_2, \dots, x_m])$ . Now we call  $L$  and  $y_i$  **vector functions**. Hence, by EQUATION(1a) and (1b), we have

$$\begin{array}{l} \text{Chain rule for} \\ \text{vector function} \end{array} : \quad \frac{\partial L}{\partial x_k} = \frac{\partial L}{\partial \mathbf{y}} \cdot \frac{\partial \mathbf{y}}{\partial x_k} \quad (4)$$

### 3 Derivatives for Matrix in Neural Network

In neural network, we assume the prior layer, say,  $(w-1)$ -th layer, outputs  $m$  data  $\mathbf{X}$ , namely  $\mathbf{X} = [\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(m)}]$  where  $\mathbf{x}^{(i)} = [x_1^{(i)}, x_2^{(i)}, \dots, x_n^{(i)}]^T$  ( $i = 1, 2, \dots, m$ , and  $x_1^{(i)}, x_2^{(i)}, \dots, x_n^{(i)}$  are  $m$  attributes or features of  $\mathbf{x}^{(i)}$ ), to the  $w$ -th layer. The  $w$ -th layer takes linear operation on  $\mathbf{X}$  and outputs  $\mathbf{Z}$ :

$$\mathbf{Z} = \mathbf{W}\mathbf{X} + \mathbf{B} \quad (5)$$

where  $\mathbf{W}$  is weight matrix and  $\mathbf{B}$  is bias matrix:

$$\mathbf{Z} = [\mathbf{z}^{(1)}, \mathbf{z}^{(2)}, \dots, \mathbf{z}^{(m)}] \quad (6a)$$

$$\mathbf{W} = \begin{bmatrix} w_{11} & w_{12} & \cdots & w_{1n} \\ w_{21} & w_{22} & \cdots & w_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ w_{p1} & w_{p2} & \cdots & w_{pn} \end{bmatrix}, \quad \mathbf{B} = \underbrace{\begin{bmatrix} b_1 & b_1 & \cdots & b_1 \\ b_2 & b_2 & \cdots & b_2 \\ \vdots & \vdots & \ddots & \vdots \\ b_p & b_p & \cdots & b_p \end{bmatrix}}_{m \text{ columns}} \quad (6b)$$

We can multiply  $\mathbf{W}$  to each entry of  $\mathbf{X}$  and therefore,

$$\mathbf{z}^{(k)} = \mathbf{W}\mathbf{x}^{(k)} + \mathbf{b} \quad (k = 1, 2, \dots, m) \quad (7)$$

where  $\mathbf{b} = [b_1, b_2, \dots, b_p]^T$ .

### 3.1 For single datum

For simplicity, let's first consider the case of only one datum, the  $k$ -th example:  $\mathbf{z}^{(k)} = \mathbf{W}\mathbf{x}^{(k)} + \mathbf{b}$ , namely

$$\begin{bmatrix} z_1^{(k)} \\ z_2^{(k)} \\ \vdots \\ z_p^{(k)} \end{bmatrix} = \begin{bmatrix} w_{11} & w_{12} & \cdots & w_{1n} \\ w_{21} & w_{22} & \cdots & w_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ w_{p1} & w_{p2} & \cdots & w_{pn} \end{bmatrix} \begin{bmatrix} x_1^{(k)} \\ \vdots \\ x_n^{(k)} \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_p \end{bmatrix} \quad (8)$$

It's obvious that  $(i = 1, 2, \dots, p)$

$$z_i^{(k)} = \sum_{j=1}^n w_{ij} x_j^{(k)} + b_i \quad (9)$$

Let  $L^{(k)} : \mathbb{R}^p \mapsto \mathbb{R}$  be  $L^{(k)} = L^{(k)}(\mathbf{z}^{(k)})$ . We will evaluate  $\frac{\partial L^{(k)}}{\partial \mathbf{W}}$ ,  $\frac{\partial L^{(k)}}{\partial \mathbf{x}^{(k)}}$  and  $\frac{\partial L^{(k)}}{\partial \mathbf{b}}$ . By EQUATION(3), we have

$$\frac{\partial L^{(k)}}{\partial \mathbf{W}} = \begin{bmatrix} \frac{\partial L^{(k)}}{\partial w_{11}} & \frac{\partial L^{(k)}}{\partial w_{21}} & \cdots & \frac{\partial L^{(k)}}{\partial w_{p1}} \\ \frac{\partial L^{(k)}}{\partial w_{12}} & \frac{\partial L^{(k)}}{\partial w_{22}} & \cdots & \frac{\partial L^{(k)}}{\partial w_{p2}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial L^{(k)}}{\partial w_{1n}} & \frac{\partial L^{(k)}}{\partial w_{2n}} & \cdots & \frac{\partial L^{(k)}}{\partial w_{pn}} \end{bmatrix} \quad (10)$$

Now let's take one entry from matrix above, say,  $\frac{\partial L^{(k)}}{\partial w_{ij}}$ . Since  $L^{(k)} = L^{(k)}(\mathbf{z}^{(k)})$  and  $\mathbf{z}^{(k)} = \mathbf{z}^{(k)}(\mathbf{W})$ , by EQUATION(4), we have that

$$\frac{\partial L^{(k)}}{\partial w_{ij}} = \frac{\partial L^{(k)}}{\partial \mathbf{z}^{(k)}} \frac{\partial \mathbf{z}^{(k)}}{\partial w_{ij}} \quad (11)$$

where

$$\frac{\partial L^{(k)}}{\partial \mathbf{z}^{(k)}} = \left[ \frac{\partial L^{(k)}}{\partial z_1^{(k)}}, \frac{\partial L^{(k)}}{\partial z_2^{(k)}}, \dots, \frac{\partial L^{(k)}}{\partial z_p^{(k)}} \right] \quad (12a)$$

$$\frac{\partial \mathbf{z}^{(k)}}{\partial w_{ij}} = \left[ \frac{\partial z_1^{(k)}}{\partial w_{ij}}, \frac{\partial z_2^{(k)}}{\partial w_{ij}}, \dots, \frac{\partial z_p^{(k)}}{\partial w_{ij}} \right]^T \quad (12b)$$

Using EQUATION(9), there is, for  $h = 1, 2, \dots, p$

$$\begin{aligned} \frac{\partial z_h^{(k)}}{\partial w_{ij}} &= \frac{\partial}{\partial w_{ij}} \left( \sum_{t=1}^n w_{ht} x_t^{(k)} + b_h \right) \\ &= \sum_{t=1}^n \frac{\partial w_{ht}}{\partial w_{ij}} x_t^{(k)} \\ &= \sum_{t=1}^n \delta_{ij,ht} x_t^{(k)} \\ &= \delta_{i,h} x_j^{(k)} \end{aligned} \quad (13)$$

This infers that the  $i$ -th entry of  $\frac{\partial \mathbf{z}^{(k)}}{\partial w_{ij}}$  is  $x_j^{(k)}$  while the others are 0, and so we conclude that

$$\frac{\partial L^{(k)}}{\partial w_{ij}} = \frac{\partial L^{(k)}}{\partial z_i^{(k)}} x_j^{(k)} \quad (14)$$

Hence,

$$\begin{aligned} \frac{\partial L^{(k)}}{\partial \mathbf{W}} &= \begin{bmatrix} \frac{\partial L^{(k)}}{\partial z_1^{(k)}} x_1^{(k)} & \frac{\partial L^{(k)}}{\partial z_2^{(k)}} x_1^{(k)} & \dots & \frac{\partial L^{(k)}}{\partial z_p^{(k)}} x_1^{(k)} \\ \frac{\partial L^{(k)}}{\partial z_1^{(k)}} x_2^{(k)} & \frac{\partial L^{(k)}}{\partial z_2^{(k)}} x_2^{(k)} & \dots & \frac{\partial L^{(k)}}{\partial z_p^{(k)}} x_2^{(k)} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial L^{(k)}}{\partial z_1^{(k)}} x_n^{(k)} & \frac{\partial L^{(k)}}{\partial z_2^{(k)}} x_n^{(k)} & \dots & \frac{\partial L^{(k)}}{\partial z_p^{(k)}} x_n^{(k)} \end{bmatrix} \\ &= \begin{bmatrix} x_1^{(k)} \\ x_2^{(k)} \\ \vdots \\ x_n^{(k)} \end{bmatrix} \begin{bmatrix} \frac{\partial L^{(k)}}{\partial z_1^{(k)}} & \frac{\partial L^{(k)}}{\partial z_2^{(k)}} & \dots & \frac{\partial L^{(k)}}{\partial z_p^{(k)}} \end{bmatrix} \\ &= \mathbf{x}^{(k)} \frac{\partial L^{(k)}}{\partial \mathbf{z}^{(k)}} \end{aligned} \quad (15)$$

Since the derivative of scalar value w.r.t. a matrix is still a matrix whose size is as the same as that of the transpose of the origin one, it's sometimes very convenient in machine learning to have the derivative matrix size invariant. Here, we define

$$\nabla_{\mathbf{W}} L^{(k)} = \left( \frac{\partial L^{(k)}}{\partial \mathbf{W}} \right)^T = \left( \frac{\partial L^{(k)}}{\partial \mathbf{z}^{(k)}} \right)^T \left( \mathbf{x}^{(k)} \right)^T \quad (16)$$

Now let's calculate  $\frac{\partial L^{(k)}}{\partial \mathbf{b}}$ . Similarly, for  $i = 1, 2, \dots, p$ , we have

$$\begin{aligned}
\frac{\partial L^{(k)}}{\partial b_i} &= \frac{\partial L^{(k)}}{\partial \mathbf{z}^{(k)}} \frac{\partial \mathbf{z}^{(k)}}{\partial b_i} \\
&= \frac{\partial L^{(k)}}{\partial \mathbf{z}^{(k)}} \left[ \frac{\partial z_h^{(k)}}{\partial b_i} \right]_{(h=1,2,\dots,p)}^{p \times 1} \\
&= \frac{\partial L^{(k)}}{\partial \mathbf{z}^{(k)}} \left[ \frac{\partial}{\partial b_i} \left( \sum_{t=1}^n w_{ht} x_t^{(k)} + b_h \right) \right]_{(h=1,2,\dots,p)}^{p \times 1} \\
&= \frac{\partial L^{(k)}}{\partial \mathbf{z}^{(k)}} [\delta_{h,i}]_{(h=1,2,\dots,p)}^{p \times 1} \\
&= \frac{\partial L^{(k)}}{\partial z_i^{(k)}} \quad \left( \begin{array}{l} \text{The } i\text{-th component of vector is 1, while} \\ \text{the others are 0's.} \end{array} \right)
\end{aligned} \tag{17}$$

and thus

$$\frac{\partial L^{(k)}}{\partial \mathbf{b}} = \left[ \frac{\partial L^{(k)}}{\partial z_1^{(k)}}, \frac{\partial L^{(k)}}{\partial z_2^{(k)}}, \dots, \frac{\partial L^{(k)}}{\partial z_p^{(k)}} \right] = \frac{\partial L^{(k)}}{\partial \mathbf{z}^{(k)}} \tag{18}$$

we define

$$\nabla_{\mathbf{b}} L^{(k)} = \left( \frac{\partial L^{(k)}}{\partial \mathbf{b}} \right)^T = \left( \frac{\partial L^{(k)}}{\partial \mathbf{z}^{(k)}} \right)^T \tag{19}$$

Finally, let's look at  $\frac{\partial L^{(k)}}{\partial \mathbf{x}^{(k)}}$ . Firstly, according to chain rule, we have

$$\frac{\partial L^{(k)}}{\partial \mathbf{x}^{(k)}} = \frac{\partial L^{(k)}}{\partial \mathbf{z}^{(k)}} \frac{\partial \mathbf{z}^{(k)}}{\partial \mathbf{x}^{(k)}} \tag{20}$$

by EQUATION(2),

$$\begin{aligned}
\frac{\partial \mathbf{z}^{(k)}}{\partial \mathbf{x}^{(k)}} &= \left[ \frac{\partial z_i^{(k)}}{\partial x_j^{(k)}} \right]_{ij}^{p \times n} \\
&= \left[ \frac{\partial}{\partial x_j^{(k)}} \left( \sum_{t=1}^n w_{it} x_t^{(k)} + b_i \right) \right]_{ij}^{p \times n} \\
&= \left[ \sum_{t=1}^n w_{it} \delta_{t,j} \right]_{ij}^{p \times n} \\
&= [w_{ij}]_{ij}^{p \times n} = \mathbf{W}
\end{aligned} \tag{21}$$

What a elegant result it is! We can rewrite this formula as

$$\frac{\partial \mathbf{z}^{(k)}}{\partial \mathbf{x}^{(k)}} = \frac{\partial (\mathbf{W} \mathbf{x}^{(k)} + \mathbf{b})}{\partial \mathbf{x}^{(k)}} = \mathbf{W} \tag{22}$$

For digression, this is far analgous to single-variable derivative of linear function  $y = kx + b$ , which is  $dy/dx = k$ .

Now back to our topic, we gat

$$\frac{\partial L^{(k)}}{\partial \mathbf{x}^{(k)}} = \frac{\partial L^{(k)}}{\partial \mathbf{z}^{(k)}} \mathbf{W} \quad (23)$$

and define

$$\nabla_{\mathbf{x}^{(k)}} L^{(k)} = \left( \frac{\partial L^{(k)}}{\partial \mathbf{x}^{(k)}} \right)^T = \mathbf{W}^T \left( \frac{\partial L^{(k)}}{\partial \mathbf{z}^{(k)}} \right)^T \quad (24)$$

### 3.2 For $m$ data

At present, let's consider a batch of samples together assembled in matrix form

$$\mathbf{Z} = \mathbf{W}\mathbf{X} + \mathbf{B} \quad (25)$$

Suppose there are  $m$  samples, the  $m$  columns of  $\mathbf{Z}$ , and let  $L = f(\mathbf{L}) = f(L^{(1)}, L^{(2)}, \dots, L^{(m)})$  where  $L^{(k)} = L^{(k)}(\mathbf{z}^{(k)})$ . Up to now, we just put a simple  $L$  definition into consideration, which we define as

$$L = \frac{1}{m} \sum_{k=1}^m L^{(k)} \quad (26)$$

Evaluate  $\frac{\partial L}{\partial \mathbf{W}}$ ,  $\frac{\partial L}{\partial \mathbf{X}}$  respectively.

By EQUATION(15), we have

$$\frac{\partial L}{\partial \mathbf{W}} = \frac{1}{m} \sum_{k=1}^m \frac{\partial L^{(k)}}{\partial \mathbf{W}} = \frac{1}{m} \sum_{k=1}^m \mathbf{x}^{(k)} \frac{\partial L^{(k)}}{\partial \mathbf{z}^{(k)}} \quad (27)$$

Beacuse the  $(i, j)$ -entry of  $\mathbf{x}^{(k)} \frac{\partial L^{(k)}}{\partial \mathbf{z}^{(k)}}$  is  $\frac{\partial L^{(k)}}{\partial z_j^{(k)}} x_i^{(k)}$ , the sum of these  $m$  matrices turns out to be

$$\begin{aligned} \sum_{k=1}^m \mathbf{x}^{(k)} \frac{\partial L^{(k)}}{\partial \mathbf{z}^{(k)}} &= \left[ \sum_{k=1}^m \frac{\partial L^{(k)}}{\partial z_j^{(k)}} x_i^{(k)} \right]_{ij}^{n \times p} \\ &= \left[ x_i^{(k)} \right]_{ik}^{n \times m} \left[ \frac{\partial L^{(k)}}{\partial z_j^{(k)}} \right]_{kj}^{m \times p} \\ &= \mathbf{X} \frac{\partial L}{\partial \mathbf{Z}} \end{aligned} \quad (28)$$

so we have

$$\nabla_{\mathbf{W}} L = \left( \frac{\partial L}{\partial \mathbf{W}} \right)^T = \frac{1}{m} \left( \frac{\partial L}{\partial \mathbf{Z}} \right)^T \mathbf{X}^T = \frac{1}{m} (\nabla_{\mathbf{Z}} L) \mathbf{X}^T \quad (29)$$

Similarly, we have

$$\nabla_{\mathbf{X}} L = \frac{1}{m} \mathbf{W}^T (\nabla_{\mathbf{Z}} L) \quad (30)$$